

United States
Department of
Agriculture

National
Agricultural
Statistics
Service

Research and
Applications
Division

SRB Research Report
Number SRB-89-12

October 1989

A COMPUTER ALGORITHM FOR MARKOV CHAIN FORECASTS OF COTTON OBJECTIVE YIELD

James H. Matis
Charles R. Perry
Donald E. Boudreaux

A COMPUTER ALGORITHM FOR MARKOV CHAIN FORECASTS OF COTTON OBJECTIVE YIELD

by James H. Matis*, Charles R. Perry and Donald E. Boudreaux*, National Agricultural Statistics Service, U.S. Department of Agriculture, Washington, D.C. 20250, October 1989, Research Report No. SRB 89-12.

Abstract

This paper documents the computer algorithm developed by Matis, Perry, Boudreaux and Aune (1989) in evaluating a revised Markov chain procedure for forecasting final cotton objective yield. The algorithm was executed in three programs. Each program's function is summarized by a statement of purpose, a procedure outline, and a few comments. The complete code along with a detailed annotation is provided for each program.

* Dr. Matis and Mr. Boudreaux are with Texas A&M University, Department of Statistics, College Station, Texas 77843.

Keywords: Markov Chain, Cross Validation, Forecast Errors, Objective Yield.

This paper was prepared for limited distribution to the research community outside the U.S. Department of Agriculture (USDA). The views expressed herein are not necessarily those of the National Agricultural Statistics Service (NASS) or USDA. The use of company names in this publication is for identification only and does not imply endorsement by the Department of Agriculture.

Acknowledgements

The authors express their appreciation to Ben Klugh and George Hanuschak for their helpful suggestions and moral support during this project. We express our gratitude to our colleagues Bill Donaldson, Barry Ford, and Phil Kott for their thoughtful reviews of this report. However, we bear full responsibility for any errors.

Contents

Summary	1
Introduction	1
Computer Algorithm Overview	1
Program: XDAT	2
Program: XRSQ	7
Program: YRXX	8
References	18

Summary

Existing computer programs were revised and expanded to accomplish two objectives. The first objective was to implement a new procedure for defining the categorical Markov transition states. The second objective was to automate the variable selection procedure so that the predictor variables are determined solely on the basis of the statistical evidence from the data at the time of forecast.

The revised algorithm was divided into five basic steps: read and edit the data, select the predictor variables, create the new categorical states, calculate the transition matrices, and generate the forecasts and estimate the forecast errors. The five steps were executed in three computer programs. Each program is described by a statement of purpose, a procedure outline, and a few comments. The complete Statistical Analysis System (SAS) code is given for each program along with a detailed annotation.

Introduction

This paper provides documentation for the algorithm and programs developed by Matis, Perry, Boudreaux and Aune (1989) in evaluating a revised Markov chain procedure for forecasting final cotton objective yield. Computer programs existing from earlier related cooperative research between the National Agricultural Statistics Service (NASS) and Texas A and M University (TAMU) (Matis, *et. al.*, 1985, 1989) were revised and expanded to accomplish two objectives. The first was to implement the new procedure for defining states. The second was to automate the program so that the primary and secondary predictor variables for each forecast were determined without human intervention. Previously these predictor variables were selected on the basis of expert judgment of the user. The revised procedures select the variables based solely on the statistical evidence available at the time of each forecast.

Computer Algorithm Overview

The process of constructing and evaluating Markov forecasts was accomplished in a five steps.

1. Read and edit the data.
2. Select independent variables.
3. Create new categorical variables based upon the selected variables and user defined number of breaks.
4. Calculate the Markov transition matrices.
5. Generate the estimated forecasts and check the accuracy of the simulation.

These steps were executed within three SAS programs – XDAT (step 1.), XRSQ (step 2.), and YRXX (steps 3., 4., 5.). Each program is described by its basic purpose, a overview of the procedure, a detailed annotated program outline, and a few comments. The complete SAS code is provided for each program. However, it should be pointed out that the programs were written for researching the feasibility of utilizing the Markov chain forecast procedure on USDA/NASS data. Thus, the programs are not necessarily optimally

coded nor are they intended for "production" use. Furthermore, the programs were written in an older mainframe version of SAS, therefore they would have to be updated for use with PC-SAS or newer versions of mainframe SAS. This is especially evident with the replacement of PROC MATRIX with PROC IML.

Program: XDAT

Purpose:

This program edits the USDA/NASS cotton objective yield SAS data sets (from tape) for consistency among the values of the various variables, create new variables, and develop a data structure appropriate for Markov analysis.

Procedure:

The original tapes contained a series of SAS data sets each representing a single year. These yearly data sets were read and processed to identify and categorize the monthly sequence of information associated with each individual plot. A unique identification number, ID, and month variable, MONTH, were created. Some data items were combined, edited, and used to create new variables. Other data sets were created and merged to obtain an overall cumulative yield variable. Certain independent variables were then ranked and a sequence of data steps were executed to organize and rename the independent variables by month (8,9, and 10). Finally, these monthly data sets were recombined and a SAS data set was created on a mainframe disk pack.

Program Outline:

Lines 001–009 JCL card for SAS batch processing.

Lines 010–013 titles to document SAS output.

Lines 014–026 macro _MKMO determines state and year it will also begin the process of creating the variable MONTH.

Lines 027–035 read 1980 data and create MONTH and ID.

Lines 036–044 read 1981 data and create MONTH and ID.

Lines 045–053 read 1982 data and create MONTH and ID.

Lines 054–062 read 1983 data and create MONTH and ID.

Lines 063–071 read 1984 data and create MONTH and ID.

Lines 072–080 read 1985 data and create MONTH and ID.

Lines 081–089 read 1986 data and create MONTH and ID.

Lines 090–137 combine 80 to 86 data sets, edit, and create new variables (note cumulative variables BURR, OPEN, and YLD).

Lines 138–142 segment out the final cumulative yield.

Lines 143–151 recombine the final yield for each observation and select independent variables to keep.

Lines 152–160 rank the independent variables.

Lines 161–170 segment out month 8 and rename variables.

Lines 171–180 segment out month 9 and rename variables.

Lines 181–190 segment out month 10 and rename variables.

Lines 191-197 recombine the data sets for months 8,9, and 10 then save the result to mainframe disk space.

Line 198 printout 20 observations to check data.

Comments:

The creation of the variables MONTH and ID was necessary because of the way the data was accumulated and structured for analysis. The macro (Lines 014-026) and data set steps (Lines 027-089) used in creating these variables are not efficient procedures. They were employed to overcome a SAS problem utilizing a "by" statement within the macro language. Attention should be given to the use of SAS's automatic data set naming convention within this particular set on previous page of data steps and procedures (Lines 031, 040, 049, 058, 067, 076, 084).

SAS Code follows:

```
001 //XDAT JOB (C774,2C,2,25,JM), 'BOUDREAUX',  
002 // MSGCLASS=Z,MSGLEVEL=(0,0)  
003 // EXEC SAS,OPTIONS=MACRO,REGION=4096K  
004 //TAPE DD DSN=USR.E413.JM.USDA1.INDEX1.YANG,DISP=SHR  
005 //DSK DD DSN=USR.E413.JM.STATE48,  
006 // DISP=(NEW,CATLG,DELETE),  
007 // UNIT=SYSDA,  
008 // SPACE=(TRK,(100,20))  
009 //SYSIN DD *  
  
010 TITLE1 ' MARKOV PROJECT      ' ;  
011 TITLE2 ' DATA FILTER PROC (NEW) ' ;  
012 TITLE3 ' AUG 88, MATIS & BOUDREAUX ' ;  
013 TITLE4 ' STATE=48(ALL)   YRS 81-84 ' ;  
  
014 %MACRO _MKMO ;  
015   IF STATE = 48 ;  
016   IF YEAR = 0 THEN DELETE ;  
017   IF YEAR = 5 THEN DELETE ;  
018   IF YEAR = 6 THEN DELETE ;  
019   N = _N_ ;  
020   PROC SORT ;  
021     BY SAMPLE ;  
022   PROC RANK ;  
023     BY SAMPLE ;  
024     VAR N ;  
025     RANKS RNK ;  
026 %MEND _MKMO ;  
  
027 DATA X80 ;  
028   SET TAPE.CTWK80 ;  
029   %_MKMO ;  
030 DATA Y80 ;  
031   SET DATA1 ;
```

```
032      ID = 80000 + SAMPLE ;
033      MONTH = RNK + 7 ;
034      IF MONTH GE 14 THEN
035          DELETE ;

036  DATA X81 ;
037      SET TAPE.CTWK81 ;
038      %_MKMO ;
039  DATA Y81 ;
040      SET DATA2 ;
041      ID = 81000 + SAMPLE ;
042      MONTH = RNK + 7 ;
043      IF MONTH GE 14 THEN
044          DELETE ;

045  DATA X82 ;
046      SET TAPE.CTWK82 ;
047      %_MKMO ;
048  DATA Y82 ;
049      SET DATA3 ;
050      ID = 82000 + SAMPLE ;
051      MONTH = RNK + 7 ;
052      IF MONTH GE 14 THEN
053          DELETE ;

054  DATA X83 ;
055      SET TAPE.CTWK83 ;
056      %_MKMO ;
057  DATA Y83 ;
058      SET DATA4 ;
059      ID = 83000 + SAMPLE ;
060      MONTH = RNK + 7 ;
061      IF MONTH GE 14 THEN
062          DELETE ;

063  DATA X84 ;
064      SET TAPE.CTWK84 ;
065      %_MKMO ;
066  DATA Y84 ;
067      SET DATA5 ;
068      ID = 84000 + SAMPLE ;
069      MONTH = RNK + 7 ;
070      IF MONTH GE 14 THEN
071          DELETE ;

072  DATA X85 ;
073      SET TAPE.CTWK85 ;
074      %_MKMO ;
075  DATA Y85 ;
076      SET DATA6 ;
```

```

077      ID = 85000 + SAMPLE ;
078      MONTH = RNK + 7 ;
079      IF MONTH GE 14 THEN
080          DELETE ;

081 DATA X86 ;
082     SET TAPE.CTWK86 ;
083     %_MKMO ;
084 DATA Y86 ;
085     SET DATA7 ;
086     ID = 86000 + SAMPLE ;
087     MONTH = RNK + 7 ;
088     IF MONTH GE 14 THEN
089         DELETE ;

090 DATA DSO ;
091     SET Y80 Y81 Y82 Y83 Y84 Y85 Y86 ;
092     IF (C509 = 0) THEN
093         CONFAC - 0 ;
094     ELSE
095         CONFAC = C510 / C509 ;
096     CURRWT = 1.0526 * CONFAC *
097     (C316+C317+C325+C326+C327+C336+C337+C345+C346+C347) ;
098     ROWSP = (C303+C304) / 8 ;
099     IF (MONTH = 8) THEN DO ;
100         BURR = C312+C322+C332+C342 ;
101         OPEN = C313+C314+C323+C324+C333+C334+C343+C344 ;
102         IF (OPEN > 0 AND CURRWT = 0) THEN
103             CURRWT = . ;
104             CUMWT = CURRWT ;
105             IF (ROWSP = 0) THEN
106                 ROWSP = 3.225 ;
107             END ;
108             ELSE DO ;
109                 BURR = C312+C322+C332+C342+ BURR ;
110                 OPEN = C313+C314+C323+C324+C333+C334+C343+C344+OPEN ;
111                 IF (OPEN > 0 AND CURRWT = 0) THEN
112                     CURRWT = . ;
113                     CUMWT = CUMWT+CURRWT ;
114                     IF (ROWSP = 0 OR ROWSP = .) THEN
115                         ROWSP = ROWXX ;
116                     END ;
117                     ROWXX = ROWSP ;
118                     IF (OPEN = 0) THEN
119                         Z21 = 0 ;
120                     ELSE
121                         Z21 = CUMWT / OPEN ;
122                         UNOPB = C319+C329+C339+C349 ;
123                         PARTB = C318+C328+C338+C348 ;
124                         LB = BURR+OPEN +PARTB +UNOPB ;

```

```

125      P3 = C350+C365 ;
126      P = C311 + C321 + C331 + C341 ;
127      IF (P = 0) THEN
128          MATUR = 0 ;
129      ELSE
130          MATUR = LB / P ;
131          X1 = (0.870 * LB) + (0.867 * (C366+C367)) ;
132          X2 = (6.667 * C368) ;
133          X3 = 6.667*(C364+C374) ;
134          MO_YLD = 2.401 * 0.368 * CUMWT / ROWSP ;
135          IF (MO_YLD NE .) THEN
136              YLD = MO_YLD ;
137          RETAIN BURR OPEN CUMWT ROWXX YLD ;

138  DATA DSF ;
139      SET DSO ;
140      IF MONTH = 13 ;
141      YIELD = YLD ;
142      KEEP ID YIELD ;

143  DATA DS2 ;
144      MERGE DSO DSF ; BY ID ;
145      IF (YIELD = 0 OR YIELD = .) THEN
146          DELETE ;
147      KEEP BURR OPEN UNOPB PARTB LB
148          X1 X2 X3 P P3
149          MATUR CONFACT CURRWT CUMWT ROWSP
150          Z21 MO_YLD C380 MONTH STATE
151          YEAR YLD YIELD ID ;

152  PROC RANK DATA=DS2 OUT=DSRK ;
153      VAR BURR OPEN UNOPB PARTB LB
154          X1 X2 X3 P P3
155          MATUR CONFACT CURRWT CUMWT ROWSP
156          Z21 C380 YIELD ;
157      RANKS RBURR ROPEN RUNOPB RPARTB RLB
158          RX1 RX2 RX3 RP RP3
159          RMATUR RCONFAC RCURRWT RCUMWT RROWS
160          RZ21 RC380 RYIELD ;

161  DATA M08 ;
162      SET DSRK ;
163      IF MONTH = 8 ;
164      RBUR_8 = RBURR ; ROPEN_8 = ROPEN ; RUNO_8 = RUNOPB ;
165      RPRT_8 = RPARTB ; RLB_8 = RLB ; RX1_8 = RX1 ;
166      RX2_8 = RX2 ; RX3_8 = RX3 ; RP_8 = RP ;
167      RP3_8 = RP3 ; RMAT_8 = RMATUR ; RCON_8 = RCONFAC ;
168      RCUR_8 = RCURRWT ; RCUM_8 = RCUMWT ; RROW_8 = RROWS ;
169      RZ21_8 = RZ21 ; RC380_8 = RC380 ;
170      KEEP ID RBUR_8 -- RC380_8 ;

```

```

171 DATA M09 ;
172   SET DSRK ;
173   IF MONTH = 9 ;
174     RBUR_9 - RBURR      ; ROPEN_9 - ROPEN      ; RUNO_9 - RUNOPB ;
175     RPRT_9 - RPARTB    ; RLB_9 - RLB        ; RX1_9 - RX1 ;
176     RX2_9 - RX2       ; RX3_9 - RX3        ; RP_9 - RP ;
177     RP3_9 - RP3       ; RMAT_9 - RMATUR    ; RCON_9 - RCONFAC ;
178     RCUR_9 - RCURRWT  ; RCUM_9 - RCUMWT    ; RROW_9 - RROWSP ;
179     RZ21_9 - RZ21     ; RC380_9 - RC380    ;
180   KEEP ID RBUR_9 -- RC380_9 ;

181 DATA M10 ;
182   SET DSRK ;
183   IF MONTH = 10 ;
184     RBUR_10- RBURR      ; ROPEN_10- ROPEN      ; RUNO_10 - RUNOPB ;
185     RPRT_10- RPARTB    ; RLB_10 - RLB        ; RX1_10 - RX1 ;
186     RX2_10 - RX2       ; RX3_10 - RX3        ; RP_10 - RP ;
187     RP3_10 - RP3       ; RMAT_10 - RMATUR    ; RCON_10 - RCONFAC ;
188     RCUR_10- RCURRWT  ; RCUM_10 - RCUMWT    ; RROW_10 - RROWSP ;
189     RZ21_10- RZ21     ; RC380_10- RC380    ;
190   KEEP YEAR YIELD ID RBUR_10 -- RC380_10 ;

191 PROC SORT DATA=M08 ; BY ID ;
192 PROC SORT DATA=M09 ; BY ID ;
193 PROC SORT DATA=M10 ; BY ID ;
194 DATA DSK.STATE482 ;
195   MERGE M08 M09 M10 ; BY ID ;
196   IF (YIELD = .) THEN
197     DELETE ;
198 PROC PRINT DATA = DSK.STATE482 (OBS = 20) ;

```

Program: XRSQ

Purpose:

This program uses the data created in XDAT to select the two best variables for predicting yield.

Procedure:

Given the nature of the study, one specific year was always excluded to allow the remaining data to be used to simulate a prediction. Then a PROC RSQUARE was run for each of the monthly time frames (8,9, and 10). The resulting set of "best" two variable models were used later in the Markov process.

Program Outline:

Lines 001-004 are JCL cards for batch processing.

Lines 005-007 are titles that will document the SAS output.

Lines 008-010 segment the year 1981 out of the data (as an example).

Lines 011-013 select the model for month 8.
Lines 014-016 select the model for month 9.
Lines 017-019 select the model for month 10.

Comments:

This process provides for an objective methodology of variable selection. However, it does assume that an "acceptable" set of independent variables are utilized.

SAS Code follows:

```
001 //RQ81 JOB (C774,2C,1,25,JM), 'BOUDREAUX'  
002 // EXEC SAS,REGION=1024K  
003 //DSK DD DSN=USR.E413.JM.STATE48,DISP=SHR  
004 //SYSIN DD *  
  
005 TITLE1 ' Variable Selection      ' ;  
006 TITLE2 ' STATE = 48(ALL)        ' ;  
007 TITLE3 ' YEARS = 82,83,84       ' ;  
  
008 DATA DS_ONE ;  
009     SET DSK.STATE482 ;  
010     IF (YEAR NE 1) ;  
  
011 PROC RSQUARE DATA=DS_ONE START=2 STOP=2 SELECT=1 CP ;  
012     MODEL YIELD = ROPEN_8    RUNO_8    RLB_8    RCUR_8  
013                 RX3_8      RP_8      RMAT_8   RZ21_8   RCUM_8 ;  
  
014 PROC RSQUARE DATA=DS_ONE START=2 STOP=2 SELECT=1 CP ;  
015     MODEL YIELD = ROPEN_9    RUNO_9    RLB_9    RCUR_9  
016                 RX3_9      RP_9      RMAT_9   RZ21_9   RCUM_9 ;  
  
017 PROC RSQUARE DATA=DS_ONE START=2 STOP=2 SELECT=1 CP ;  
018     MODEL YIELD = ROPEN_10   RUNO_10   RLB_10   RCUR_10  
019                 RX3_10    RP_10    RMAT_10  RZ21_10  RCUM_10 ;
```

Program: YRXX

Purpose:

This program takes the data from the program XDAT along with the selected variables from the program XRSQ and performs several Markov analysis forecast simulations.

Procedure:

The data was broken into two data sets one (ds_one) with the data to provide the forecasts and the other (ds_two) with the "actual" yield values, for later cross-validation analysis, (Efron 1982). Both of these data sets had new variables added to them based upon the selected variables, a set of user defined categories, and the breakpoints of the variable categories for the forecasts data set. These new categorical variables were then formed into Markov transition matrices and used to forecast the "actual" yields.

Program Outline:

Lines 001-005 JCL card for SAS batch processing.

Lines 006-009 titles to document SAS output.

Lines 010-015 read disk data and segment it into a forecast data set (ds_one) and an actual data set (ds_two).

Lines 016-079 macro _CAT1 : used to create a one level categorical break for a single variable.

Lines 080-187 macro _CAT2 : used to create a two level nested categorical break.

Lines 188-191 these lines are where the user inputs the selected variables and the number of splits per variable that are to be used.

Lines 192-203 the index variables are incremented by one for future readability.

Lines 204-207 this PROC FREQ provides the cell counts for the transition matrices.

Lines 208-216 univariate statistics are calculated for the final yield categories.

Lines 217-220 the statistics are printed and merged back into the original data set.

Lines 221-234 macro _MTX : creates the transition matrices.

Lines 235-297 macro _CMX : calculates the cumulative distribution and other statistics for the forecasted data.

Lines 298-315 macro _ESX : determines the estimate to be used (mean value or median), the actual, and the residual values. Univariate statistics and plots are then produced for the forecast simulation.

Lines 316-323 these lines are the code and macro combinations necessary to "run" the markov simulation for months 8, 9, 10 → forecast.

Lines 324-330 simulation for months 9, 10 → forecast.

Lines 331-336 simulation for month 10 → forecast.

Comments:

This SAS program was developed to satisfy a specific research need. Thus, it is not optimally coded and should not be used for production work. There are parts of the code that may be highly dependent on the version of SAS used. These include Lines 221-234 where the order of the information read from the PROC FREQ may vary and the use of PROC MATRIX in Lines 016-079, Lines 080-187, or in Lines 316-336 which has been replaced by IML. Also in some cases the number of categories asked for by the user in Lines 188-191 may not be supported by the data. When this occurs the PROC MATRIX procedures will fail to form the transition matrices. This can easily be corrected by choosing a smaller number of classifications and resubmitting the program.

SAS Code follows:

```
001 //YR81 JOB (C774,2C,1,25,JM),'BOUDREAUX',
002 //      MSGCLASS=Z,MSGLEVEL=(0,0)
003 // EXEC SAS,OPTIONS=MACRO,REGION=2048K
004 //DSK   DD DSN=USR.E413.JM.STATE48,DISP=OLD
005 //SYSIN DD *
```

```

006 TITLE1 ' MARKOV ANALYSIS      ' ;
007 TITLE2 ' XCAT & XMAT PROCEDURES   ' ;
008 TITLE3 ' AUG 88, MATIS & BOUDREAUX ' ;
009 TITLE4 ' STATE=48(ALL) MODEL YR=81  ' ;

010 DATA DS_ONE ;
011     SET DSK.STATE48 ;
012     IF (YEAR NE 1) ;
013 DATA DS_TWO ;
014     SET DSK.STATE48 ;
015     IF (YEAR EQ 1) ;

016 %MACRO _CAT1(V1,N1,BK) ;
017 /* CLASSIFY DATA IN A SECOND DATA SET BY
018 /* ONE VARIABLE IN DATA SET ONE .
019 /*
020 /* V1      : FIRST VARIABLE
021 /* N1      : NO. OF CATEGORIES TO BREAK VAR1
022 /* BK      : RESULTING CATEGORICAL VARIABLE
023 */

024 DATA DS_ONE ;
025     SET DS_ONE ;
026     IF (&V1 = .) THEN
027         DELETE ;

028 DATA DS_TWO ;
029     SET DS_TWO ;
030     IF (&V1 = .) THEN
031         DELETE ;

032 PROC RANK DATA=DS_ONE OUT=RS1 GROUPS=&N1 ;
033     VAR &V1 ;
034     RANKS BREAK1 ;
035 PROC SORT DATA=RS1 ;
036     BY BREAK1 ;
037 PROC UNIVARIATE NOPRINT DATA=RS1 ;
038     VAR &V1 ;
039     OUTPUT OUT=US1 MAX=MAX1 ;
040     BY BREAK1 ;

041 PROC MATRIX ;
042     FETCH MAX1 DATA=US1 (KEEP=BREAK1 MAX1) ;

043     SPLIT1 = J.(1,&N1,0) ;
044     SPLIT1(1,1) = 0 ;
045     DO I = 2 TO &N1 ;
046         SPLIT1(1,I) = MAX1(I-1,2) ;
047     END ;

```

```

048      PRINT SPLIT1 ;

049      FETCH NEW DATA=DS_TWO (KEEP=&V1) ;
050      NR = NROW (NEW) ;
051      YY = J.(NR,2,0) ;
052      DO I = 1 TO NR ;
053          DO J = 1 TO &N1 ;
054              IF (J LT &N1) THEN DO ;
055                  IF (NEW(I,1) GT SPLIT1(1,J)) AND
056                      (NEW(I,1) LE SPLIT1(1,J+1)) THEN
057                      YY(I,1) = J ;
058                  END ;
059              ELSE IF (J EQ &N1) THEN DO ;
060                  IF (NEW(I,1) GT SPLIT1(1,J)) THEN
061                      YY(I,1) = J ;
062                  END ;
063              ELSE
064                  YY(I,1) = . ;
065          END ;
066      END ;
067      OUTPUT YY OUT=YY ;

068      DATA YY ;
069          SET YY ;
070          &BK = COL1 - 1 ;
071          DROP COL1 ;

072      DATA DS_ONE ;
073          SET RS1 ;
074          &BK = BREAK1 ;
075          DROP BREAK1 ;

076      DATA DS_TWO ;
077          MERGE DS_TWO YY ;
078          DROP ROW ;
079  %MEND _CAT1 ;

080  %MACRO _CAT2(V1,N1,V2,N2,BK) ;
081      /* CLASSIFY DATA IN A SECOND DATA SET BY
082      /* TWO VARIABLES (NESTED) IN DATA SET ONE .
083      /*
084      /* V1      : FIRST VARIABLE
085      /* N1      : NO. OF CATEGORIES TO BREAK VAR1
086      /* V2      : SECOND VARIABLE
087      /* N2      : NO. OF CATEGORIES TO BREAK VAR2
088      /* BK      : RESULTING CATEGORICAL VARIABLE
089      */

```

```

090      DATA DS_ONE ;
091          SET DS_ONE ;
092          IF (&V1 = . OR &V2 = .) THEN
093              DELETE ;

094      DATA DS_TWO ;
095          SET DS_TWO ;
096          IF (&V1 = . OR &V2 = .) THEN
097              DELETE ;

098      PROC RANK DATA=DS_ONE OUT=RS1 GROUPS=&N1 ;
099          VAR &V1 ;
100          RANKS BREAK1 ;
101      PROC SORT DATA=RS1 ;
102          BY BREAK1 ;
103      PROC RANK DATA=RS1 OUT=RS2 GROUPS=&N2 ;
104          VAR &V2 ;
105          RANKS BREAK2 ;
106          BY BREAK1 ;
107      PROC SORT DATA=RS2 ;
108          BY BREAK1 BREAK2 ;
109      PROC UNIVARIATE NOPRINT DATA=RS2 ;
110          VAR &V1 ;
111          OUTPUT OUT=US1 MAX=MAX1 ;
112          BY BREAK1 ;
113      PROC UNIVARIATE NOPRINT DATA=RS2 ;
114          VAR &V2 ;
115          OUTPUT OUT=US2 MAX=MAX2 ;
116          BY BREAK1 BREAK2 ;

117      PROC MATRIX ;
118          FETCH MAX1 DATA=US1 (KEEP=BREAK1           MAX1) ;
119          FETCH MAX2 DATA=US2 (KEEP=BREAK1 BREAK2 MAX2) ;
120          NR = NROW (MAX2) ;

121          SPLIT1 = J.(1,&N1,0) ;
122          SPLIT1(1,1) = 0 ;
123          DO I = 2 TO &N1 ;           ;
124              SPLIT1(1,I) = MAX1(I-1,2) ;
125          END ;
126          PRINT SPLIT1 ;

127          SPLIT2 = J.(&N1,&N2,0) ;
128          DO I = 1 TO NR ;
129              R = MAX2(I,1) + 1 ;
130              C = MAX2(I,2) + 1 ;
131              SPLIT2 (R,C) = MAX2(I,3) ;
132          END ;

```

```

133      DO I = 1 TO &N1 ;
134          DO J = &N2 TO 2 BY -1 ;
135              SPLIT2 (I,J) = SPLIT2 (I,J-1) ;
136          END ;
137          SPLIT2 (I,1) = 0.0 ;
138      END ;
139      PRINT SPLIT2 ;

140      FETCH NEW DATA=DS_TWO (KEEP=&V1 &V2) ;
141      NR = NROW (NEW) ;
142      YY = J.(NR,2,0) ;
143      DO I = 1 TO NR ;
144          DO J = 1 TO &N1 ;
145              IF (J LT &N1) THEN DO ;
146                  IF (NEW(I,1) GT SPLIT1(1,J)) AND
147                      (NEW(I,1) LE SPLIT1(1,J+1)) THEN
148                      YY(I,1) = J ;
149                  END ;
150              ELSE IF (J EQ &N1) THEN DO ;
151                  IF (NEW(I,1) GT SPLIT1(1,J)) THEN
152                      YY(I,1) = J ;
153                  END ;
154              ELSE
155                  YY(I,1) = . ;
156          END ;
157          W = J.(1,1,0) ;
158          W = YY(I,1) ;
159          DO J = 1 TO &N2 ;
160              IF (J LT &N2) THEN DO ;
161                  IF (NEW(I,2) GT SPLIT2(W,J)) AND
162                      (NEW(I,2) LE SPLIT2(W,J+1)) THEN
163                      YY(I,2) = J ;
164                  END ;
165              ELSE IF (J EQ &N2) THEN DO ;
166                  IF (NEW(I,2) GT SPLIT2(W,J)) THEN
167                      YY(I,2) = J ;
168                  END ;
169              ELSE
170                  YY(I,2) = . ;
171          END ;
172      END ;
173      OUTPUT YY OUT=YY ;

174      DATA YY ;
175          SET YY ;
176          COL1 = COL1 - 1 ;
177          COL2 = COL2 - 1 ;
178          &BK = &N2*COL1 + COL2 ;
179          DROP COL1 COL2 ;

```

```

180      DATA DS_ONE ;
181          SET RS2 ;
182          &BK = &N2*BREAK1 + BREAK2 ;
183          DROP BREAK1 BREAK2 ;

184      DATA DS_TWO ;
185          MERGE DS_TWO YY ;
186          DROP ROW ;
187  %MEND _CAT2 ;

188  %_CAT2(RX3_8    ,4 ,RLB_8    ,2 ,INX1) ;
189  %_CAT2(RLB_9    ,4 ,RX3_9    ,2 ,INX2) ;
190  %_CAT2(RUNO_10   ,8 ,RCUM_10   ,2 ,INX3) ;
191  %_CAT1(YIELD,50,SY1 ) ;

192  DATA ONE ;
193      SET DS_ONE ;
194      INX1 = INX1 + 1 ;
195      INX2 = INX2 + 1 ;
196      INX3 = INX3 + 1 ;
197      SY1 = SY1 + 1 ;
198  DATA TWO ;
199      SET DS_TWO ;
200      INX1 = INX1 + 1 ;
201      INX2 = INX2 + 1 ;
202      INX3 = INX3 + 1 ;
203      SY1 = SY1 + 1 ;

204  PROC FREQ DATA=ONE ;
205      TABLES INX1*INX2 / NOPRINT OUT=T12 ;
206      TABLES INX2*INX3 / NOPRINT OUT=T23 ;
207      TABLES INX3*SY1 / NOPRINT OUT=T3Y ;

208  PROC SORT DATA=ONE ; BY SY1 ;
209  PROC UNIVARIATE NOPRINT DATA=ONE ; BY SY1 ;
210      VAR YIELD ;
211      OUTPUT OUT=UNI
212          MEAN   = MNY1
213          MEDIAN = MDY1
214          Q1     = Q1
215          Q3     = Q3
216          VAR    = VAR ;

217  PROC PRINT DATA=UNI ;
218  PROC SORT DATA=UNI ; BY SY1 ;
219  DATA ONE ;
220      MERGE ONE UNI ; BY SY1 ;

```

```

221  %MACRO _MTX (FRQ,AMX,A_ROW,A_COL) ;
222      /* READ PROC FREQ DATA INTO MATRICIES
223      /*
224      /* FRQ      : PROC FREQ OUTPUT DATASETS
225      /* AMX      : RESULTING MATRIX
226      */
227      FETCH &FRQ DATA=&FRQ ;
228      &AMX = J.(&A_ROW,&A_COL,0) ;
229      DO I = 1 TO NROW(&FRQ) ;
230          &AMX(&FRQ(I,1),&FRQ(I,2)) = &FRQ(I,3) ;
231      END ;
232      /* SCALE THE ROWS OF MATRIX : SUM = 1.0 */
233      &AMX = (1 #/ (DIAG(&AMX(,+)))) * &AMX ;
234  %MEND _MTX ;

235  %MACRO _CMX (AMX,CMX,STA,INX) ;
236      /* GENERATE STATISTICS
237      /*
238      /* AMX      : ORIGINAL MATRIX
239      /* CMX      : CUM DIST OF ORIGINAL MATRIX
240      /* STA      : Q1,Q3,MEDIAN,MEAN,VAR OF AMX
241      /* INX      : INDEX NUMBER, MUST CORRESPOND TO AMX
242      */
243      FETCH UNI DATA=UNI (KEEP=MNY1) ;
244      NC = NCOL(&AMX) ;
245      NR = NROW(&AMX) ;
246      &CMX = J.(NR,NC+1,0) ;
247      &CMX = &AMX || J.(NR,1,0) ;
248      &STA = J.(NR,6,0) ;
249      DO I = 1 TO NR ;
250          CUM = 0 ;
251          VAR = 0 ;
252          MEAN = 0 ;
253          DO J = 1 TO NC ;
254              MEAN = MEAN + (UNI(J,1)*&AMX(I,J)) ;
255              VAR = VAR + ((UNI(J,1)**2)*&AMX(I,J)) ;
256              CUM = CUM + &AMX(I,J) ;
257              &CMX(I,J) = CUM ;
258              IF (CUM LE 0.25) THEN &STA(I,1) = J ;
259              IF (CUM LE 0.50) THEN &STA(I,2) = J ;
260              IF (CUM LE 0.75) THEN &STA(I,3) = J ;
261          END ;
262          &CMX(I,NC+1) = I ;
263          IF (&STA(I,1) = NC) THEN
264              GOTO LX ;
265          /* *** Q1      */
266          &STA(I,1) = UNI(&STA(I,1),1)
267          + (UNI(&STA(I,1)+1,1)
268          - UNI(&STA(I,1),1))

```

```

269      * ((0.25 - &CMX(I,&STA(I,1)))
270      #/ (&CMX(I,&STA(I,1)+1)
271      - &CMX(I,&STA(I,1)))) ;
272 /* *** MEDIAN */
273 &STA(I,2) = UNI(&STA(I,2),1)
274      + (UNI(&STA(I,2)+1,1)
275      - UNI(&STA(I,2),1))
276      * ((0.50 - &CMX(I,&STA(I,2)))
277      #/ (&CMX(I,&STA(I,2)+1)
278      - &CMX(I,&STA(I,2)))) ;
279 /* *** Q3   */
280 &STA(I,3) = UNI(&STA(I,3),1)
281      + (UNI(&STA(I,3)+1,1)
282      - UNI(&STA(I,3),1))
283      * ((0.75 - &CMX(I,&STA(I,3)))
284      #/ (&CMX(I,&STA(I,3)+1)
285      - &CMX(I,&STA(I,3)))) ;
286 /* *** MEAN */
287 &STA(I,4) = MEAN ;
288 /* *** VAR */
289 &STA(I,5) = VAR - MEAN**2 ;
290 /* *** INDEX */
291 &STA(I,6) = I    ;
292 LX : END ;
293 OUTPUT &STA OUT=&STA
294      (RENAME = (COL1=Q1 COL2=MEDIAN COL3=Q3
295                  COL4=MEAN COL5=VAR COL6=&INX)) ;
296 FREE CUM MEAN VAR ;
297 %MEND _CMX ;

298 %MACRO _ESX (EST,STA,INX) ;
299 /* MERGE STATISTICS WITH NEW DATA
300 */
301 /* EST      : WHICH ESTIMATE (MEAN, MEDIAN)
302 /* STA      : DATASET OF MATRIX STATISTICS
303 /* INX      : WHICH INDEX VARIABLE TO USE
304 */
305 PROC PRINT DATA=&STA ;
306 PROC SORT DATA=&STA ; BY &INX ;
307 PROC SORT DATA=TWO ; BY &INX ;
308 DATA TWO ;
309      MERGE TWO &STA ; BY &INX ;
310      EST = &EST      ; /* ESTIMATES YIELD */
311      ACT = YIELD    ; /* ACTUAL YIELD */
312      RES = EST - ACT ; /* RESIDUAL */
313      PROC UNIVARIATE PLOT ;
314          VAR EST ACT RES ;
315 %MEND _ESX ;

```

```
316 PROC MATRIX ;
317   %_MTX (T12,A1,8,8)
318   %_MTX (T23,A2,8,16)
319   %_MTX (T3Y,A3,16,50)
320   EST = A1*A2*A3 ;
321   %_CMX (EST,CDF,STATS,INX1)
322   %_ESX (MEAN,STATS,INX1)
323   TITLE5 ' MONTHS 8,9,10      ';

324 PROC MATRIX ;
325   %_MTX (T23,A2,8,16)
326   %_MTX (T3Y,A3,16,50)
327   EST = A2*A3 ;
328   %_CMX (EST,CDF,STATS,INX2)
329   %_ESX (MEAN,STATS,INX2)
330   TITLE5 ' MONTHS 9,10      ';

331 PROC MATRIX ;
332   %_MTX (T3Y,A3,16,50)
333   EST = A3 ;
334   %_CMX (EST,CDF,STATS,INX3)
335   %_ESX (MEAN,STATS,INX3)
336   TITLE5 ' MONTHS 10      ';
```

References

- Efron, B. (1982). The Jackknife, the Bootstrap and Other Resampling Plans. Society for Industrial and Applied Mathematics, Philadelphia, PA.
- Matis, J. H., Saito, T., Grant, W. E. Iwig, W. C., and Richie, J. T. (1985). A Markov chain approach to crop yield forecasting. *Agricultural Systems* 18:171-187.
- Matis, J. H., Birkett, T., and Boudreux, D. (1989). An application of the Markov Chain approach to forecasting cotton yield from surveys. *Agricultural Systems* 29:357-370.
- Matis, J.H., Perry, C.R., Boudreux, D.E., and D.J. Aune (1989). *Markov Chain Forecasts of Cotton Objective Yield*, National Agricultural Statistics Service, U.S. Department of Agriculture, Washington, D.C. 20250, Research Report No. SRB 89-11.

U. S. GOVERNMENT PRINTING OFFICE:1989-251-433; 07704/NA35